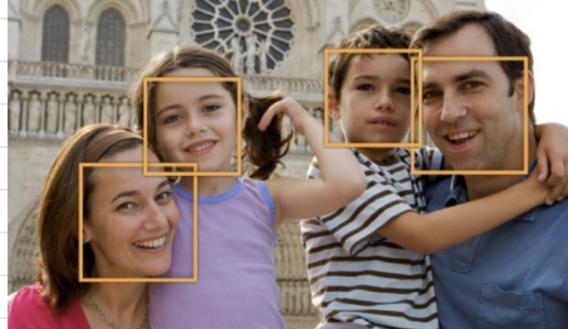


LECTURE 2 : FACE DETECTION

#1



Face detection is part of object detection/recognition, which most of the techniques are based on machine learning.



Machine Learning's Common framework

Training stage
(supervised learning)

Input : data + labels

Output : model parameters

Testing stage :

Input : query data + model parameters

Output : predicted labels

Data = images

→
transform to

Features = Important factors
of images



Face detection technique : Viola - Jones Face Detection.

Basic Components : (1) Features
(2) Classification

Feature Extraction : Image Convolution

» 1D Convolution :

Analog operation : $(f * g)(x) = \int f(u) g(x-u) du$

Discrete operation : $(f * g)(x) = \sum_{u=-\infty}^{+\infty} f(u) g(x-u)$

e.g. :

$$f = \text{the 1D kernel} = \begin{array}{|c|c|c|} \hline u: & f_1 & 0 & -1 \\ \hline 1 & 2 & 3 \\ \hline \end{array}$$

$$g = \text{the input signal} = \begin{array}{|c|c|c|c|c|c|} \hline x: & 1 & 2 & 3 & 4 & 5 \\ \hline 5 & 0 & 7 & 0 & \dots \\ \hline \end{array}$$

$$(f * g)(x=2) = (3 \times 7) + (2 \times 10) + (1 \times 5) = 46$$

$$(f * g)(x=3) = (3 \times 0) + (7 \times 2) + (1 \times 10) = 24$$

$$h(x) = (f * g)(x) = \text{the output} = \boxed{0 \ 46 \ 24 \ \dots}$$

» 2D Convolution :

$$(f * g)(x, y) = \sum_{v=-\infty}^{+\infty} \sum_{u=-\infty}^{+\infty} f(u, v) g(x-u, y-v)$$

$f(u, v)$ = the 2D kernel

e.g. :

$$k = \begin{array}{|c|c|c|} \hline -1 & 0 & +1 \\ \hline -2 & 0 & +2 \\ \hline -1 & 0 & +1 \\ \hline \end{array}$$

$\downarrow v \quad \rightarrow u$

$$I = \text{Input image} : (k * I)(x, y) = \sum_v \sum_u k(u, v) I(x-u, y-v)$$

See on the internet {

- sobel operator
- Gaussian kernel / filter
- Laplacian kernel / filter

[2] Feature Extraction: Haar-like Features

#2

① Three types of masks :

② Convolution with the masks:

③ Descriptor:

Type 1:

$$I = \begin{array}{|c|c|c|c|c|} \hline \textcircled{1} & \textcircled{2} & & & \\ \hline & & & & \\ \hline \end{array}$$

Type 2:

e.g.: I 

Type 3:

Assuming pixel ① has intensity 100 and pixel ② 150, then:

$$F = (100 \times -1) + (150 \times +1)$$

= 50

Another example:

$$I = \begin{array}{|c|c|c|c|} \hline 1 & 2 & 3 & 4 \\ \hline 5 & 6 & 7 & 8 \\ \hline \end{array} \quad * \quad \begin{array}{|c|c|c|} \hline \text{black} & \text{black} & \text{white} \\ \hline \text{black} & \text{black} & \text{white} \\ \hline \end{array}$$

$$F = [\textcircled{3} + \textcircled{4} + \textcircled{7} + \textcircled{8}] - [\textcircled{1} + \textcircled{2} + \textcircled{5} + \textcircled{6}]$$

For an image of size 24×24 , the convolution of the image with all types of the masks & all possible sizes will generate 160,000 values of F (= features)

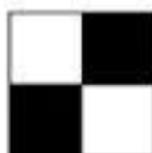
This vector is called
a descriptor



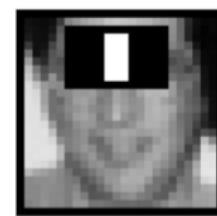
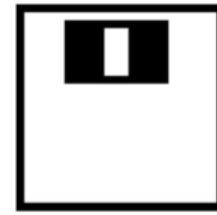
(a) Edge Features



(b) Line Features

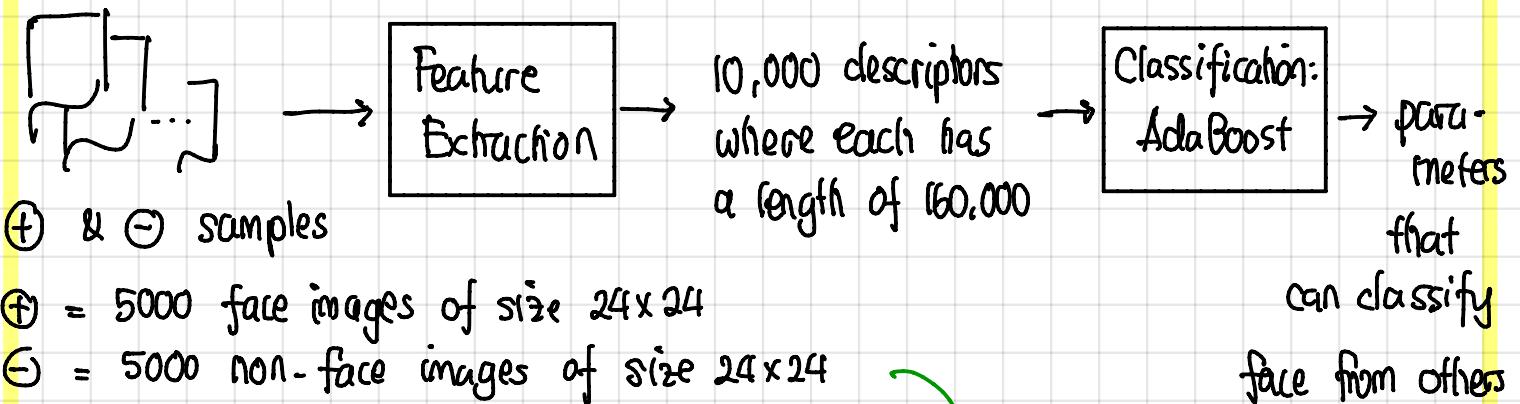


(c) Four-rectangle features



[3] Training Pipeline

#3



The positive (+) samples:



The code for generating 160,000 features from a single 24x24 image:

```
const int frameSize = 24;
const int features = 5;
// All five feature types:
const int feature[features][2] = {{2,1}, {1,2}, {3,1}, {1,3}, {2,2}};

int count = 0;
// Each feature:
for (int i = 0; i < features; i++) {
    int sizeX = feature[i][0];
    int sizeY = feature[i][1];
    // Each position:
    for (int x = 0; x <= frameSize-sizeX; x++) {
        for (int y = 0; y <= frameSize-sizeY; y++) {
            // Each size fitting within the frameSize:
            for (int width = sizeX; width <= frameSize-x; width+=sizeX) {
                for (int height = sizeY; height <= frameSize-y; height+=sizeY) {
                    count++;
                }
            }
        }
    }
}
```

[4] AdaBoost

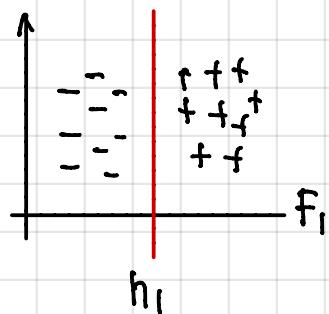
#4

Goal : to find functions that separate positive & negative samples

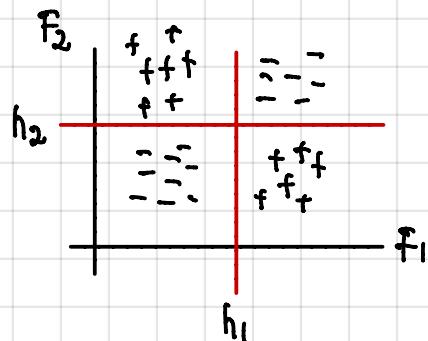
Basic idea: to have a strong classifier (or function) composed of a number of weak classifiers.

Examples :

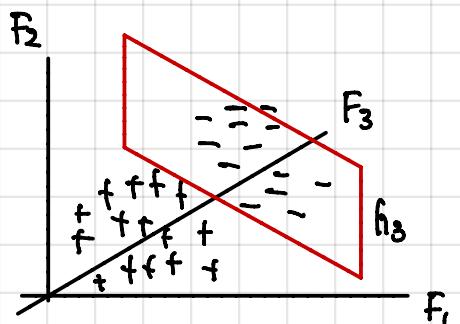
1D data



2D data

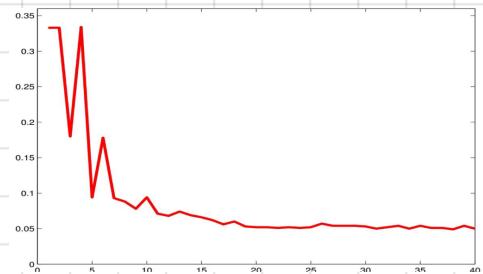
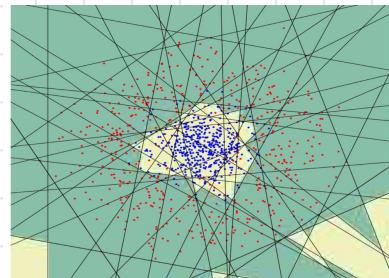
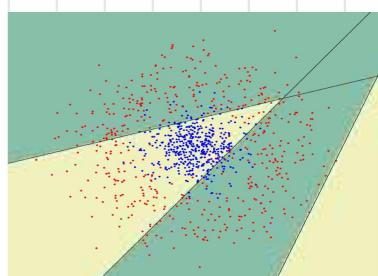
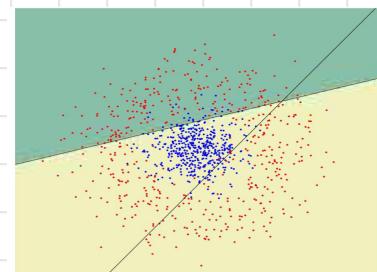
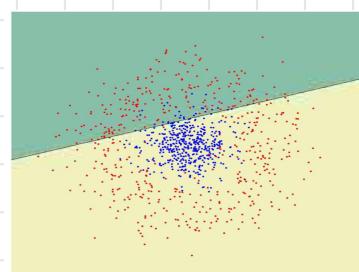
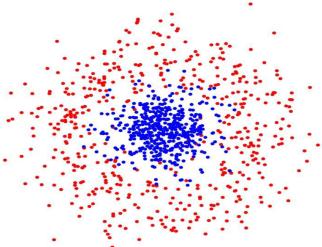


3D data



1. h_1, h_2, h_3, \dots are called weak classifiers that can form a strong classifier.
2. In our case, we have 160,000 features, meaning we have 160,000 dimensions of space.

Illustration of the general AdaBoost algorithm:



[5] AdaBoost: Weights

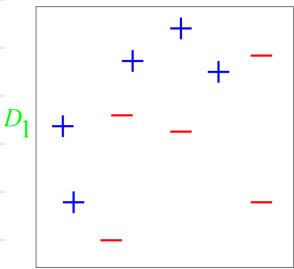
#5

AdaBoost has important properties:

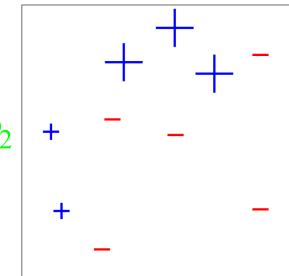
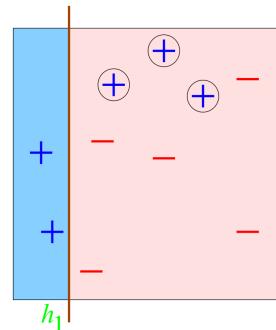
1. Iterative in nature
2. In each iteration, it generates a weak classifier
3. In the next iteration, it focuses on samples that have more weights.

Q: What does it mean? What are the weights?

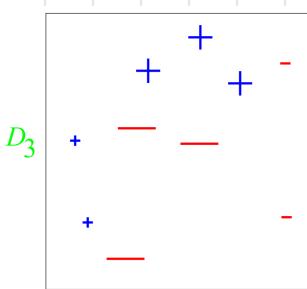
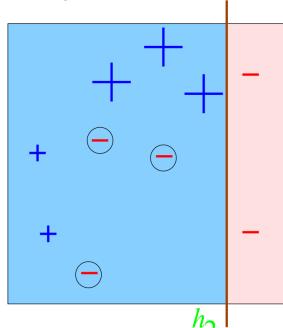
A: Consider the following example:



Round 1:

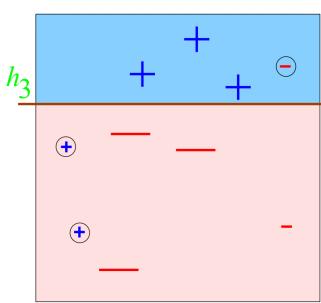


Round 2:

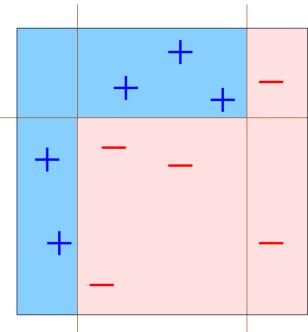


wrongly classified samples are given more weights, so that in the next round (Round 2), they will get more attention.

Round 3:



Final hypothesis:



→ a strong classifier, which is obtained by combining the weak classifiers (h₁, h₂, h₃)

[6] Training Algorithm based on AdaBoost

#6

1. Input: descriptor & labels: $(x_1, y_1), \dots, (x_N, y_N)$

x_i = a descriptor (160,000 length features) ; $N = 10,000$

$y_i = \begin{cases} 0 & : \text{negative sample} \\ 1 & : \text{positive sample} \end{cases}$; $i = \text{the index of samples or descriptors}$

2. Initialize weights: $w_{1,i} = \begin{cases} 1/m & \ominus \\ 1/l & \oplus \end{cases}$; $m, l = \text{the numbers of negative \& positive samples, respectively.}$

3. For $t = 1, \dots, T$:

($T = \text{the number of the weak classifiers / dimensions of data}$)

(1) Normalize the weights: $w_{t,i} = \frac{w_{t,i}}{\sum_{j=1}^N w_{t,j}}$

(2) Select the best weak classifier: the one that minimizes the error:

$$\epsilon_t = \min_{\{\theta\}} \sum_{i=1}^N w_{t,i} |h(x_i, \theta) - y_i|$$

! \rightarrow every zero if x_i is correctly classified

range: $0 \sim 1$ all functions/ sample
(due to $w_{t,i}$) parameters of possible weak classifiers

(3) Define: $h_t(\bar{x}) = h(\bar{x}, \theta_t)$

either 0 or 1 (represents the weak classifier) \hookrightarrow the best weak classifier that minimizes ϵ_t

(4) Update the weights: $w_{t+1,i} = w_{t,i} \beta_t^{(1-\alpha_i)}$

where: $\alpha_i = \begin{cases} 0 & \text{if } x_i \text{ is correctly classified} \\ 1 & \text{else} \end{cases}$

1. $x_i = 1$ (incorrect classification) $\rightarrow w_{t+1}$ doesn't change

2. $\alpha_i = 0$ (correct classification) & ϵ_t is small (< 0.5), w_{t+1} becomes smaller

3. $x_i = 0$ & $\epsilon_t > 0.5$, w_{t+1} becomes larger as $\beta_t \gg 1$ \rightarrow to penalize &

give x_i the same as case #1.

$$\beta_t = \log \frac{1}{\beta_t}$$

4. Strong classifier:

$$C(x_{\text{test}}) = \begin{cases} 1 & \text{if } \sum_{t=1}^T \beta_t h_t(x_{\text{test}}) \geq \frac{1}{2} \sum_{t=1}^T \beta_t \\ 0 & \text{else} \end{cases}$$

$\frac{1}{2} \sum_{t=1}^T \beta_t$ give more weight to a better h_t

[7] Testing: Simple Version

Input: a single image of size 24×24

Output: a label $\begin{cases} 1 & \text{if face} \\ 0 & \text{else} \end{cases}$

Algorithm:

- ① Calculate the Haar-like descriptor from the input image
- ② Use AdaBoost & the trained parameters to decide / classify.

But, this is not what we want.

What we want: Given an image in any size, to know the locations of the faces (if any) and to draw bounding boxes surrounding those faces.



Naive Solution

Input: any image of any size (gray)

Algorithm

Problem: SLOW

- ① Generating subwindows
- ② Checking using the strong classifier

Integral Image

↓

Cascade

- (1) Scan each pixel in the input image, and take a patch / subwindow with size of 24×24 pixels.
- (2) Using the strong classifier, check if any of the subwindows is face.
- (3) Increase the size of the subwindows i.e. 30×30 pixels, scan the entire pixels, and check again if any subwindows are faces.
- (4) Keep increasing the size of the subwindows till a certain size (= max. face size)

[9] Integral Image

1. Integral image formulation:

$$ii(x, y) = ii(x-1, y) + s(x, y)$$

$$s(x, y) = s(x, y-1) + i(x, y)$$

Original input image

$$s(x, -1) = 0$$

$$ii(-1, y) = 0$$

Example:

Input image, i :

0	1
5	2
3	6



ii : the integral image:

5
8
16

$$ii(0,1) = ii(-1,1) + s(0,1)$$

= 0

$$s(0,1) = s(0,0) + i(0,1)$$

= 3

$$s(0,0) = s(0,-1) + i(0,0) = 5$$

= 0

$$ii(0,1) = 0 + 5 + 3 = 8$$

$$ii(1,1) = ii(0,1) + s(1,1)$$

= 8

$$s(1,1) = s(0,1) + i(1,1)$$

= 6

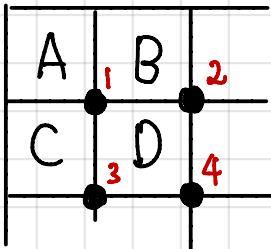
$$s(0,1) = s(1,-1) + i(1,0)$$

= 0

= 2

$$ii(1,1) = 8 + 6 + 2 = 16$$

2. Computing a block:



- 1 = A
- 2 = A + B
- 3 = A + C
- 4 = A + B + C + D

Given the values of 1, 2, 3, 4. What is the value of block D?

$$D = 4 - (A+B+C) = 4 - (2+3) = 4 - (2 + (3-A))$$

$$= 4 - (2+3-1) = (4+1) - (2+3)$$

Example: i (input)

5	2	5	2
3	6	3	6
5	2	5	2
3	6	3	6

39 16

ii (integral image)

5	7	12	14
8	16	24	32
13	23	36	46
16	32	48	64

Total sum indicated by red points:

$$D_{red} = (4+1) - (2+3) \\ = (64+16) - (32+32) = 16$$

Blue points:

$$D_{blue} = (64+5) - (14+16) = 39$$

[10] Testing: Using Integral Image

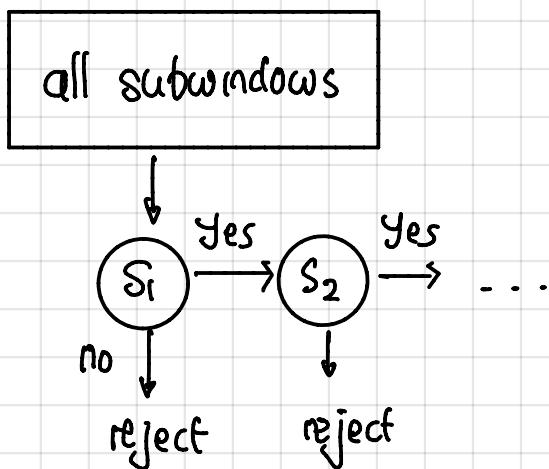
#9

Input: any gray image of any size

- Algorithm:
- (1) Compute the integral image
 - (2) Generate square windows starting from 24×24 pixels for all parts of the integral images
 - (3) Compute the descriptors using the standard mask sizes.
 - (4) Increase the sizes of the masks to capture larger square windows and compute the descriptor.
- 

This will be similar to increasing the size of subwindows on the input image, but faster.

[11] Cascade of Classifiers



S_i = a stage composed of a few weak classifiers.

1. We want to reject soon, since non-face subwindows appear more frequently than face subwindows.
2. The number of the stages is determined by training the cascade.

Additional Notes:

$$1. \quad \beta_t = \frac{\varepsilon_t}{1 - \varepsilon_t} \quad (\text{see page } 46)$$

2. A strong classifier is defined as:

$$C(x_{\text{test}}) = \begin{cases} 1 & \sum_t \beta_t h_t(x_{\text{test}}) \geq \frac{1}{2} \sum_t \beta_t \\ 0 & \text{else} \end{cases}$$

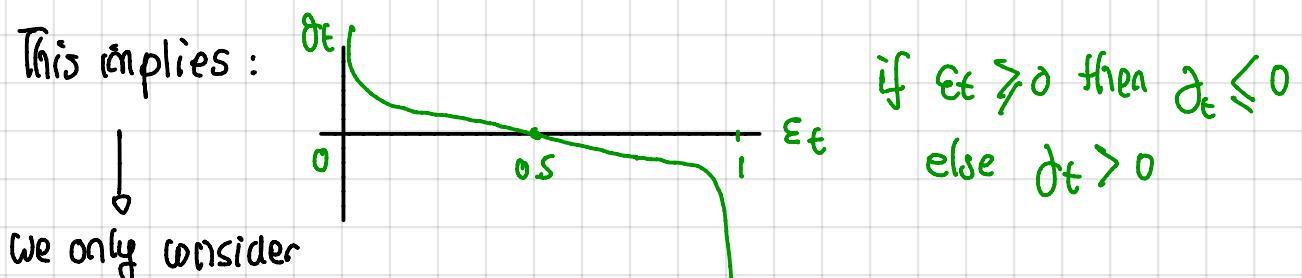
(1) Assuming $\beta_t = 1$ for every t , then the condition is:

$$\sum_t^T h_t(x_{\text{test}}) \geq \frac{1}{2} T ; \quad T = 160,000$$

This implies: For x_{test} to be labeled 1, at least it must be supported by 50% of the whole weak classifiers (up to T)

$$(2) \quad \beta_t = \log \frac{1}{\beta_t} = \log \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right) = \log \left(\frac{1}{\varepsilon_t} - 1 \right)$$

This implies:



We only consider

the weak classifiers whose error is small.

3. How to use $C(\bar{x})$, the strong classifier, in testing?

$$C(x_i, \bar{x}) = \begin{cases} 1 & \sum_t \beta_t h_t(x_i, \bar{x}) > \frac{1}{2} \sum_t \beta_t \\ 0 & \text{else} \end{cases}$$

Where $h_t(x_i, \bar{x})$ means we apply the weak classifier h_t on x_i .

$C(x_i, \bar{x})$ means we label x_i to be 1 if most of the weak classifiers that have error less than 0.5, classify x_i as a positive sample.